

Echo Chambers, Filter Bubbles, and Cyber Warfare: An Ethical and Anticipatory Ethical Analysis

Richard Wilson^{1,2} and Noah Donnelly²

¹Department of Philosophy, Towson University, Baltimore, Maryland, USA

²Computer Science and Information Sciences, Towson University, Baltimore, Maryland, USA

wilson@towson.edu

ndonnell1@students.towson.edu

Abstract: With the development of the internet in the information era and the wide access to information the internet makes available, Echo Chambers and Filter Bubbles have developed. Echo Chambers and Filter Bubbles are a consequence of Reinforcement Learning Algorithms. An Echo Chamber is an environment where people only encounter beliefs or opinions that reinforce the beliefs and opinions to which they are already committed. (Sunstein, 2017). This serves the purpose of creating a constant positive feedback loop, which continually reinforces one idea or set of ideas (Murphy, 2022). A Filter Bubble develops when recommendation algorithms feed users content based on what narratives the recommendation algorithm determines an audience wants to hear to maximize user engagement (Pariser, 2011). Echo Chambers and Filter Bubbles are relevant to Cyber and Cognitive Warfare because state actors can take advantage of the existence of these environments and use them to influence discourse and public opinion. In the context of cyber warfare attackers can use bots and fake accounts where they pretend to be citizens of the target nation to flood these Echo Chambers with narratives that align with the audience's current belief system and while also benefiting the attackers (Singer & Brooking, 2018). In the context of cyber warfare attackers can abuse Filter Bubbles by using data breaches and advertising data to infiltrate the Filter Bubbles and direct them towards the attackers desired narrative (Matz et al., 2017). These abuses rise above the level of mere internet trolling, they are intentional and targeted acts of Cyber Warfare aimed at influencing the cognitive space of a nation's population (Claverie & du Cluzel, 2022). Echo Chambers and Filter Bubbles are constructed to accomplish a strategic objective, in this case they are aimed at influencing the decision making and opinions of an audience to influence elections, policy, protests, or overall public sentiment. These strategic objectives are accomplished by using non-kinetic weapons to destabilize society, diminish decision making capabilities, and erode trust in democratic institutions. This type of cyber-attack was made apparent in the COVID-19 disinformation schemes when wellness communities on social media platforms were flooded with anti-vaccine and medicine narratives leading to distrust in the medical system, and the politicians who promoted them (Dawson & Innes, 2019). This paper identifies the technical, ethical, and anticipated ethical issues of Filter Bubbles and Echo Chambers and proposes a technical and policy framework to classify and prevent these acts of Cyber Warfare.

Keywords: Echo chambers, Filter bubbles, Cognitive warfare, Information operations, Algorithms, Micro-targeting, Anticipatory ethics, Epistemic paternalism, OCEAN model

1. Introduction

In the domain of modern conflict, perhaps the most potent form of warfare is the cognitive domain. Social media platforms and content recommendation algorithms have led to the development of two distinct yet related phenomena: Echo Chambers and Filter Bubbles. Echo Chambers are online areas where a person only encounters a point of view or opinion that mirrors their own opinions. Filter Bubbles endorse single track mindsets created by search algorithms calculating what the user wants to hear or see based on information collected about them such as: location, past-click behavior, and search history (Pariser, 2011). Initially echo chambers and filter bubbles were viewed as sociological side effects of the information era, but in truth they have now been weaponized by malicious state actors as a means of cyber warfare.

The combination of echo chambers, and filter bubbles with cyber warfare attacks represents an evolution in Information Warfare (IW). Traditional cyber warfare focused on stealing data and destroying or manipulating critical infrastructure. Echo chambers now represent a new paradigm, they focus on destroying social cohesion, and the manipulation of truth making a popular consensus on social topics difficult to attain. Attackers can insert disinformation into targeted echo chambers which can cause the public sentiment about an adversary nation to change without the use of any kinetic force. This can be called Non-Kinetic Warfare which targets the ability of a nation's people to make informed and educated decisions. These tactics can destabilize societies, influence elections, and erode trust in their democratic institutions, which align perfectly with the goals of authoritarian and fascist styles regimes (Mercy Corps, 2020).

This paper argues that echo chambers and filter bubbles are being manipulated to obscure facts and truth and the diversity of opinion and they should be recognized and classified as a form of cyber warfare, which requires a new ethical framework for assessment. On the technical side we will analyze the mechanisms which are used

to manipulate algorithms and micro-target populations. On the ethical side we will investigate the effects and implications of cognitive autonomy being eroded and apply the principles of Anticipatory Ethics to identify the moral responsibilities of platform owners. By analyzing case studies which include, Cambridge Analytica, the Myanmar genocide, and the Russian use of botnets, we show that the mode of attack in the 21st century is the human mind and not just a kinetic military base.

2. Technical Issues

Recognizing the weaponization of echo chambers and filter bubbles requires the examination of the technology that allows and facilitates their existence. The technical analysis provides the foundation for ethical and anticipatory ethical analysis of echo chambers and filter bubbles. Influencing vulnerable target audiences is an intended function of the recommendation algorithms and the data structures they rely on.

2.1 The Mathematics of Radicalization

Social media platforms use complicated Machine Learning (ML) algorithms, mainly Reinforcement Learning (RL) which is intended to maximize user engagement, to target users. The algorithms that are employed use a tactic called Reward Function, where success is measured by time-on-site or interaction (clicks, likes, shares). The algorithm observes the user's state (current mood, history), performs the action of recommending content, and receives a reward in the form of user engagement. Algorithms have learned that high-arousal emotions such as fear, outrage and validation generate the highest engagement rewards (Vosoughi et al., 2018). This creates a Positive Feedback Loop for target audiences. The algorithm feeds users increasingly extreme content to achieve higher engagement rewards, which in turn pushes the user deeper into the filter bubble, which can lead to radicalization. Radicalization can be exploited through injecting content that is mathematically optimized to trigger the highest engagement rewards, which can at the same time be radical disinformation. Botnets can artificially inflate engagement metrics of extremely divisive content, effectively gamifying the Reinforcement Learning (RL) model, which tricks the RL model into promoting disinformation to millions of real users. The algorithm is unable to recognize truth, so it becomes an unwitting technological accomplice in the radicalization of target audiences

2.2 Psychometric Micro-Targeting and the OCEAN Model

Algorithmic attacks using echo chambers, and filter bubbles are extremely precise due to the large amount of personal data on social media platforms, typically gained through data brokers or data breaches. The large amounts of data gathered about individuals allow for Psychometric Micro-targeting. This technique transcends demographics (age, location) and instead targets personality traits. This is where the OCEAN Model (McCrae & John, 1992) can be applied, and the so called Big Five personality traits can be targeted: Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism. Users of social media have a "digital exhaust" which are recorded as likes, shares and browsing history. Algorithms can predict the traits associated with the digital exhaust with a high degree of accuracy. If an attacker identifies a user who scores high on the Neuroticism scale which would make them prone to anxiety and fear, they can be targeted by disinformation employing narratives that misinform them about looming threats. An important category of psychometric micro-targeting is users with low Openness, who are targeted with more traditional propaganda. This transforms propaganda from a blunt, generalized tool into a more precise, sharp tool tailored to specific neurological vulnerabilities of a micro-targeted audience (Matz et al., 2017).

2.3 Bot-Human Hybrid Network

"Sybil Attacks" are disinformation campaigns wherein a single adversary controls multiple fake identities to manipulate the trust network metrics. (See: Gunturu, Rupesh) Modern operations use scripts from fake accounts and AI-driven personas, so when someone engages in an argument online, it may not even be with another human. Russian botnets have begun to integrate Large Language Models (LLMs) into botnets, generating contextually relevant comments that are unique, but which are created by bots. These bots can argue for prolonged periods of time within online comment sections and online forum threads which gives the appearance of a public consensus. This technique is referred to as AstroTurfing. Attackers can map the Social Graph of an echo chamber to identify Bridge Nodes which are users who have influence and who span multiple communities. By compromising or controlling Bridge Nodes attackers can infect multiple echo chambers across multiple platforms causing the misinformation presented in the echo chamber to cascade (Ferrara et al., 2016; Tomassi et al., 2024).

3. Ethical Issues

The weaponization of echo chambers and filter bubbles raises important ethical questions regarding individual autonomy and what responsibilities states have for protecting individual autonomy.

3.1 Cognitive Sovereignty and Consent

The central ethical violation in this kind of warfare is related to **what is known** as Cognitive Sovereignty, (Claverie & du Cluzel, 2022). Cognitive sovereignty is the right of an individual to control **their** mental processes and for that individual to be free opinions influenced by covert operations. When users sign up for social media, they consent to receiving targeted advertisements, but there is no clause stating that they consent to being the victim of state-sponsored psychological operations (PSYOPS). Hidden algorithms being used to curate people's **perceptions** of reality deprives them of the agency to make informed decisions. This violates Kant's 2nd Categorical Imperative which states that one must treat humanity as an end-in-itself, as opposed to a political end. Users become potential puppets for foreign governments in a conflict that they are unaware of, where their perception of reality has been altered by foreign adversaries (Claverie & du Cluzel, 2022). The question is to what degree states and social media platforms have to protect the cognitive sovereignty of individuals.

3.2 Epistemic Paternalism and the "Truth" Dilemma

Countering these attacks raises the ethical issue of Epistemic Paternalism. Platforms are effectively in a "Catch 22" dilemma where if they break the filter bubbles and downrank disinformation, they are now the arbiters of what is "true" or "false" or "good" or "bad" in the world and of what the users should see. Corporate algorithms become the chief judge of what is true and what is good information. Contrarily if the platforms do not intervene, cyber warfare is now capable of destabilizing democracy. If a social media platform does intervene, it is now acting as a paternal figure to the user, determining the information audiences are and are not allowed to see. There is not a neutral setting for these platforms, so what must be determined is what causes the least harm, and who has the moral authority to regulate these values: states, corporations, or users? (Aird, 2022).

3.3 The Erosion of Trust and Democratic Norms

Echo chambers and filter bubbles channel citizens towards different and often conflicting realities, making consensus between opposing groups nearly impossible. From a Utilitarian outcomes-based standpoint, this causes an incredible amount of harm to the stability of democratic societies, allowing democracy to erode. This channeling of citizens towards different and often conflicting realities has contributed to the erosion of trust in democratic institutions. When a population cannot agree on the most basic of facts like the effectiveness of vaccines or the legitimacy of an election, the social contract between citizens and fellow citizens and between citizens and government has broken down. Cyber warfare conducted through echo chambers and filter bubbles targets the trust of the people in the institutions of a nation, which is just as critical for a nation as its electrical grid. Unlike physical infrastructure, once this type of trust is broken, it is difficult to rebuild.

4. Case Studies

In the early stages of the development of the internet the dangers of echo chambers were previously largely theoretical, their exploitation was not yet well known, however, they recently have materialized and have shown themselves to be to date one of the most dangerous non- kinetic weapons. The following case studies will demonstrate how state and non-state actors have used the technical architecture of social media platforms to their advantage when undergoing strategic military and political operations employing echo chambers and filter bubbles.

4.1 Cambridge Analytica (2016)

Cambridge Analytica is essentially the founding father of psychometric profiling and of enabling the ability for bad actors to use psychometric profiling as a means for carrying out cognitive warfare. Cambridge Analytica collected the personal data of roughly 87 million Facebook members without the users consent through a third-party quiz app named "This Is Your Digital Life." This was not simple demographic targeting, but a much more personalized means of targeting audiences. The engineers then fed data into the OCEAN Model (Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism) to psychologically profile every voter in the United States. Through an analysis of what was dubbed a user's "digital exhaust" that included likes, shares and comments, Cambridge Analytica's algorithms were able to predict personality traits with such accuracy that they overrode the acceptance of the opinions of users own family members. The firm was able to identify "Persuadables" in swing states, who were users who scored high in Neuroticism and low in Openness. The

identified users were then fed “dark posts” which were nonpublic advertisements only visible to the target audiences. These dark posts were mathematically engineered to trigger specific anxieties and fears in audiences related to immigration or government overreach. This was an exploitation of filter bubble architecture, Cambridge Analytica helped make fictive personalized realities for millions of voters based off their deepest vulnerabilities. This was the great hack of the democratic decision-making process, voting machines were no longer the target, it was the voters mind (Cadwalladr, 2018).

4.2 Myanmar Genocide (2017)

In Myanmar, **Facebook’s engagement-based recommendation** algorithms were abused by military operatives to spark a genocide against the Rohingya Muslim minority. This is the direct, kinetic consequence of algorithmic amplification. The military of Myanmar (Tatmadaw) launched a secret operation using fake accounts and “sock puppet” pages that were disguised as entertainment or celebrity news to enter into, infiltrate and influence the digital lives of the people they were supposed to protect. Once these accounts were integrated into the user’s daily lives, they started to spread anti Rohingya Muslim propaganda labeling the Rohingya population as an existential threat to Buddhism. The technological failure lay in the hands of Facebook’s engagement ranking algorithms, which had the consequence of prioritizing hate speech since it generated intense engagement, while outrage and fear scored the highest in “time-on-site” metrics. The platform then became an enormous, radicalized echo chamber, and calls to violence were subsequently boosted by the algorithm and spread to millions of users who had no alternate news sources. This digital incitement led to physical violence, leading to the murder of at least 10,000 Rohingya people and the displacement of over 700,000 of them. This is proof of flooding filter bubbles, filled with state-actuated hate speech, which resulted in more than polarization. In the Myanmar Genocide case, the platforms code aided in a mass atrocity (Mozur, 2018).

4.3 Russian Doppelganger Operation (2022-Present): The Clone War

During February of 2022 Russia began their ongoing invasion and war against Ukraine. Since that time they had also begun a sophisticated campaign known as “Doppelganger” to destabilize the European support for Kyiv. This operation uses a “Supply Chain Attack” on truth by creating clones of legitimate Western media. Russian operatives set up domains that would try to prey on typos. They set up e. g. bild.ltd instead of bild.de, and theguardian.co.com instead of theguardian.com, and flooded these sites with fake articles (with disinformation) that mimicked the visual style including fonts and bylines of the original publications and legitimate sources. These fake articles advocated narratives that argued sanctions on Russia were destroying the European economy and that Ukrainian refugees were committing heinous crimes. If they had just deployed these websites they would gain very little traction, so in addition they had to use their botnets to flood social media sites with comments and posts that are linked to these fake sites. The sophistication of this attack lay in the ability to bypass an individual user’s skepticism filter, which creates a situation where a user who trusts The Guardian, is now likely to trust these fake sources since the disinformation presented to the audience looks like it comes from a legitimate real source. By targeting very specific filter bubbles, such as French farmers who are worried about fuel prices, or people working in the German manufacturing sector that feared layoffs, the Doppelganger operation feeds false enemy narratives straight into domestic political conversations of NATO nations, which exploits the trust of a targeted population in the free press (US Cyber Command, 2024).

4.4 COVID-19 Disinformation

During COVID-19, there was a concerted effort to radicalize new audiences by state-actors. These bad actors started targeting the “Wellness” communities filter bubbles to radicalize a new audience through the “Pastel QAnon” phenomenon. Disinformation campaigns began to target communities that were focused on yoga, organic foods, and alternative medicines, groups that already harbored a latent distrust of the pharmaceutical companies and government mandates coming from government officials. State-actors and domestic extremists used the algorithm to connect users interested in “natural immunity” with groups promoting “medical freedom,” which served in turn as the gateway to the darker conspiracy theories that talked about “The Great Reset” and QAnon messaging about globalist cabals. Russian and Chinese state actors amplified these messages to undermine Western public health responses from government officials. The “Data Void” technical exploit was where attackers generated content for unique search terms like “Plandemic” to make sure their conspiracy theories dominated search results well before they could be fact checked and before authoritative sources could respond. The wellness communities, which were previously untouched by far-right political extremism were now effectively pipelined to a point where reality became secondary to political disinformation narratives. This materialized in real world harm that included refusing vaccines, harassing medical professionals, and destabilizing public health infrastructure related to trusting healthcare information coming from government

officials. This case shows how any filter bubble can be abused, even ones focused on health and positive information, for the purposes of political extremism and destabilizing a political enemy (Claverie & du Cluzel, 2022).

5. Anticipatory Ethics

By applying Anticipatory Ethics to these case studies, what gets discovered are the moral obligations of developers of social media platforms and lawmakers related to echo chambers and filter bubbles before this next generation of warfare is fully deployed. Anticipatory ethics is a recent development in ethics concerned with examining ethical issues with technologies and technological artifacts from the research and development stage, through the introduction stage, to the stage of marketplace permeation and saturation. (Brey, 2012) Anticipatory engineering ethics and anticipatory ethical analysis attempts to identify ethical issues with Engineering and ICT technologies before and as technology develops and as they are introduced. This type of ethical analysis can also be extended to technologies related to conducting cyber warfare. (Wilson, 2021). This section will analyze these conflicts through the lens of "Moral Responsibility for Computing Artifacts" rules developed by Miller et al. (2011) and the Association for Computing Machinery (ACM) Code of Ethics.

5.1 Rule 1: The Foreseeability of Effect and ACM Principle 1.2

Rule 1 of Anticipatory Ethics employing Moral Responsibility for Computing Artifacts states that "The people who design... are morally responsible for that artifact, and for the foreseeable effects of that artifact" (Miller et al., 2011). When applied to echo chambers and filter bubbles, it is undeniably foreseeable that algorithms that are rewarded for engagement no matter the type of message, could be used by state-actors to radicalize users. The typical argument made by the companies that own the platforms who try to spin themselves as neutral town squares, is ethically invalid when the town square is designed to reward the most toxic and divisive content that they host, all in the name of "time-on-site" metrics. This design choice violates ACM Principle 1.2 "Avoid Harm" (Association for Computing Machinery, 2018). Harm is not just a metric of physical injury; it is also psychological violence. Psychological violence and the destruction of social cohesion, the incitement of violence against a specific ethnic group, and the degrading of democratic voter's abilities to make informed decisions are negative consequences. Anticipatory Ethics requires developers to move beyond reactive content moderation, which allows harm to spread, but now mandates them to implement "Red Teaming" for cognitive risks during the design phase of their technologies. This would involve simulating how bad state actors, such as the Russian Internet Research Agency, think and operate. Through this type of analysis, they could secure a new feature they are thinking about rolling out, before it is released to the public. Failing to anticipate these potential exploitations, given the amount of undeniable evidence about previous election interferences, represents a failure of their professional duty and a breach of the social contract that exists between consumer and producer as well as between producer and government.

5.2 Rule 4: The Sociotechnical Imperative and ACM Principle 3.7

Rule 4 tells us that "People who design... can do so responsibly only when they make a reasonable effort to take into account the sociotechnical systems in which the artifact is embedded" (Miller et al., 2011). The sociotechnical systems were designed for developing friendships, connections, and sharing. These platforms are currently operating in systems defined by "information warfare," "geopolitical conflict" and "hybrid threats." By not adapting their architecture to this new reality of hostile and volatile internet occurrences, there is a failure to uphold the Sociotechnical Imperative. This corresponds with ACM Principle 3.7 where professionals are required to "Recognize and take care of systems that become integrated into the infrastructure of society" (Association for Computing Machinery, 2018). Since social media has become the new public sphere where information is spread, including information about health policy, elections, and debates, developers should now carry a burden related to the duty of care. A responsible developer would implement algorithms that by design expose users to a diverse array of viewpoints, as opposed to the current algorithms that promote tunnel vision through echo chambers and that reinforce filter bubbles. When Social media companies allow their deployed algorithms to remain this way, the developers are essentially splitting the population into feudal warring tribes on a digital stage instead of a traditional battlefield. Their algorithms are inherently counterintuitive to a functioning democratic system.

5.3 Rule 2: The Post Deployment Mandate and the Problem of Epistemic Paternalism

Rule 2 tells developers that "Responsibility includes being answerable for the behaviors of the artifact... after deployment" (Miller et al., 2011). Platforms duties do not end at deployment; they have a continuing duty to

monitor their interpersonal ecosystems and moderate them. They cannot abdicate responsibility at the time of deployment. However, by strictly enforcing this mandate they introduce the issue of Epistemic Paternalism, which is the idea that platforms must intervene in the user's acquisition of knowledge "for their own good" (Aird, 2022). A counterargument posed to this could be if a platform breaks someone's filter bubble or scrubs disinformation, they act as an ultimate authority who decides what is true and what is not. This limiting of the autonomy of the user can be argued to be like how we treat children and animals, not other human beings. However, while limiting autonomy is considered unethical, an Anticipatory Ethical analysis could argue that if the user's autonomy has already been compromised by foreign state actors, engaged in attacking the platform user's psyche, the platform is morally obliged to intervene. In this context, any intervention is not to control what a user thinks, but to restore their ability to think autonomously, more specifically, to restore the user's ability to be able to think freely without interference from covert manipulation. This would morally justify aggressive actions like banning botnets, labelling state-controlled media, and friction induced sharing prompts not as acts of censorship but as a necessitated "Cognitive Defense."

6. Recommendations

To eliminate the threat of cyber warfare that utilizes echo chambers and filter bubbles, the over-reliance on passive content moderation must be abandoned, and in its place, a multi-layered defensive strategy that involves technical reforms, re-engineering of algorithms, and cognitive security protocols needs to be developed.

6.1 Diversity Aware Recommendation Systems (DARS)

The technical aspect of the target of recommendation systems used for radicalization was the "Similarity Bias" that is at the center of collaborative filtering environments. To defend against this, platforms need to implement Diversity-Aware Recommender Systems (DARS) that use Determinantal Point Processes (DPP). Basic matrix factorization models optimize only for relevance (predicting what a user is going to click). In contrast to this, a DPP-based layer would model a diverse set of items, that would ensure that the user is fed a broader semantic space as opposed to being trapped in their filter bubble. We propose that the objective function of the recommendation engine should include a "Serendipity Hyperparameter" (λ). This parameter is meant to punish semantic redundancy. If the vector embedding of the top five recommended articles are the same, it indicates identical political narratives, and the DPP layer would then introduce a repulsive force, switching redundant articles for content that is semantically distinct and highly rated by a diverse set of user clusters. This proposition would mathematically enforce the rule that people should be exposed to conflicting narratives, allowing them autonomy to choose which narratives they think are best, popping their filter bubble without requiring slow human moderation (Helberger et al., 2018).

6.2 Adversarial Training and Inverse Propensity Scoring

Adversaries attack Reinforcement Learning (RL) systems when they inject "poisoned" data into a system, which are fake likes, fake shares, and fake comments aimed at skewing the reward function. To eliminate this, platforms need to strengthen their Reinforcement Learning models with Adversarial Training. This would involve subjecting the RL model to training that would allow it to "Red Team" itself to eliminate recommendations that are likely to lead to radicalization. Platforms must also implement "Inverse Propensity Scoring" (IPS) which detects and discounts engagement metrics that are boosted by inauthentic behavior patterns, like those of bots. We recommend a "Trust Weighted Reward Function" where the platform multiplies the value of a signal such as a like or share by the users Reputational Scores which are determined by their network centrality and the age of their account. If a piece of content is being amplified by a cluster of low reputation accounts, the algorithm gives it a negative weight which suppresses the disinformation. This neutralizes Astroturfing campaigns by allowing only verified, and organic humans can be allowed to drive what is trending on a platform at any given time (Hussain et al., 2025; Platt et al., 2024).

6.3 Cryptographic Content Provenance (C2PA)

The "Doppelganger" operations and "Supply Chain Attacks" on truth must be identified and defeated. To overcome these attacks, the internet needs to adopt Cryptographic Content Provenance standards like C2PA (Coalition for Content Provenance and Authenticity). This standard would bind a cryptographically secure signature to all media files at the point of creation or publication. Respected news organizations and government agencies would be required to sign their legitimate content using their private keys. Social media platforms would then be required to implement trust indicators that act as digital watermarks and verify the chain of custody of that information. If the chain is broken, so is the digital signature, where it would lose its verified

status. So, if a Russian operative tried to clone The Guardian's website, it would be missing their signature, and the user would be alerted to the spoof. Since the user is already cognitively overloaded, the burden of verification is put upon the backs of the browsers SSL/TLS infrastructure which creates a chain of trust that state actors would be unable to create (Coalition for Content Provenance and Authenticity, 2023; Bushey, 2025).

6.4 Privacy-Preserving Algorithmic Auditing by Differential Privacy

Trust in the digital age requires independent verification, but the platforms that host their algorithms hide behind "trade secrets" and "user privacy" to prevent audits. We recommend the adoption of Differential Privacy interfaces for the purpose of auditing. This would allow independent researchers to query the social media platforms datasets to measure for the trajectory of radicalization of users without exposing any individual user's data. By adding controlled statistical noise to the query results, the platforms can prove to the regulators that their platforms are not pushing users towards extremism. This component solves the "Black Box" problem, it allows for auditing and enforcement of algorithmic liability laws without compromising their closed source software and user anonymity (O'Neil, 2016).

6.5 Cognitive Security Heuristics and Friction-by-Design

Cognitive Security (CogSec) is focused on hardening humans against emotional exploitation. Platforms must add "Friction-by-Design" heuristics. An example of this are Natural Language Processing (NLP) classifiers that are able to detect "High-Arousal/Low-Fact" linguistic patterns that are commonly found in disinformation publications. When a user attempts to share the newly flagged content without clicking on the source link, the system should throw them a "Read First" latency prompt. Platforms should also implement Inoculation Theory techniques where users are warned about disinformation tactics (e.g., "Warning: Users are currently targeting this topic with fake generated images") before they even engage with the feed. This forces users into an analytical thinking stage, breaking the outrage sharing loop (Maertens et al., 2024).

7. Future Work

Future research needs to analyze the escalation of cognitive weaponry that is driven by Generative AI. We predict the rise of "Hyper-Personalized Propaganda," where Large Language Models (LLMs) will be able to generate unique and persuasive messages for every single potential voter by scanning their real-time psychometric data, which would make "template matching" obsolete. Future work needs to explore Neuro-Symbolic AI defenses, which combine the pattern recognition abilities of neural networks and the logical reasoning of AI. Current models only understand word associations, Neuro-Symbolic AI could evaluate the factual consistency of a claim, allowing for automated fact checking. Additionally, research into Cognitive Immunology is required, specifically "AI Guardians" which are personal agents that run on a user's device to weed out manipulative content and highlight logical fallacies in real-time. This would act as a digital immune system and a potential defense for the mind of the user.

8. Conclusion

The manipulation of echo chambers and filter bubbles should be counted as a sophisticated form of cyber warfare, targeting not physical infrastructure, but now targeting the cognition of individuals instead. Through the analysis of technical mechanisms like Reinforcement Learning Loops, the OCEAN model, and botnet amplification, we demonstrate that these are not mere accidents, but engineered vulnerabilities that state-actors and radicals can use to cause harm while bypassing traditional military defenses. The case studies of Cambridge Analytica, the Myanmar genocide, and the Russian "Doppelganger" operations are proof that when content curation is left unguarded, the result is the incitement of violence with the additional result of degradation of democratic sovereignty.

This analysis also concludes that defense is in fact possible by developing strategies through the application of Anticipatory Ethics to the engineering process. Redesigning recommender systems to prioritize content diversity, enforcing Cryptographic Content Provenance, and requiring Trust-Aware Reward Functions, are all actions engineers can take to hedge these platforms against manipulation. The Sociotechnical Imperative refers to the pressing need to align technological development and social systems to achieve sustainability goals (See: , requires that developers must treat recommendation algorithms as critical infrastructure. We must now acknowledge that freedom of speech does not imply that an algorithm can recommend everything under the sun. To secure our democratic future, we need to defend the integrity of the human mind with the same rigor we apply to the defense of our military bases and outposts.

Ethics Declaration: No human participants or personally identifiable information were involved. All data sources were publicly available.

AI Tools Declaration: ChatGPT 5.1 for drafting and refinement. Human authors verified all content. Gemini 3.0 pro used for aiding in sourcing.

References

- Alabama Policy Institute. (2020, November). *Understanding the Difference Between Positive and Negative Rights*. Alabama Policy Institute.
- Association for Computing Machinery. (2018). *ACM Code of Ethics and Professional Conduct*. <https://www.acm.org/code-of-ethics>
- Aird, R. (2022). A puzzle of epistemic paternalism. *Philosophical Psychology*, 36(4), 1011–1029. <https://doi.org/10.1080/09515089.2022.2146490>
- Brey, P. A. E. (2012). Anticipatory Ethics for Emerging Technologies. *NanoEthics*, 6(1), 1–13.
- Bushey, J. (2025). Cryptographic provenance and AI-generated images. *AI Collaboratory*. https://ai-collaboratory.net/wp-content/uploads/2025/11/S13212_7356.pdf
- Cadwalladr, C. (2018). *The Great Hack: The Scandal of Cambridge Analytica*. The Observer.
- Claverie, B., & du Cluzel, F. (2022). Cognitive Warfare: The New Battlefield Exploiting Our Brains. *Polytechnique Insights*. Coalition for Content Provenance and Authenticity. (2023). *C2PA technical specification*. https://c2pa.org/specifications/specifications/1.3/specs/C2PA_Specification.html
- Collingridge, D. (1980). *The Social Control of Technology*. St. Martin's Press.
- Dawson, A., & Innes, M. (2019). How Russia's Internet Research Agency Built its Disinformation Campaign. *The Political Quarterly*, 90(2), 245–256.
- Ferrara, E., Varol, O., Davis, C., Menczer, F., & Flammini, A. (2016). The rise of social bots. *Communications of the ACM*, 59(7), 96–104.
- Gunturu, Rupesh. Survey of Sybil Attacks in Social Networks. [1504.05522](https://arxiv.org/abs/1504.05522) (accessed 2/5/2026).
- Hussain, M., Mehmood, A., Khan, M. A., Khan, R., & Lloret, J. (2025). Reputation-based leader selection consensus algorithm with rewards for blockchain technology. *Computers*, 14(1), 20. <https://doi.org/10.3390/computers14010020>
- Maertens, R., Spampatti, T., & van der Linden, S. (2024). Adding 'grit' to boost psychological inoculation against misinformation. *Journal of Experimental Psychology: Applied*. Advance online publication. <https://osf.io/dzbn7/download/?format=pdf>
- Matz, S. C., Kosinski, M., Nave, G., & Stillwell, D. J. (2017). Psychological targeting as an effective approach to digital mass persuasion. *Proceedings of the National Academy of Sciences*, 114(48), 12714–12719.
- McCrae, R. R., & John, O. P. (1992). An introduction to the five-factor model and its applications. *Journal of Personality*, 60(2), 175–215.
- Miller, K.W., et al. (2011). Moral Responsibility for Computing Artifacts: 'The Rules'. *IT Professional*, 13(3), 57–59.
- Mercy Corps. (2020). *The weaponization of social media: Landscapes assessment and case studies*. https://www.mercycorps.org/sites/default/files/2020-01/Weaponization_Social_Media_FINAL_Nov2019.pdf
- Mozur, P. (2018). *A Genocide Incited on Facebook, With Posts from Myanmar's Military*. The New York Times.
- Murphy, T. F. (2022). *Feedback loops in psychology: Positive vs. negative loops for change*. Psychology Fanatic. <https://psychologyfanatic.com/feedback-loops/>
- NATO Science & Technology Organization. (2024). *The understanding of cognitive warfare in comparative perspective*. NATO STO Meeting Proceedings. <https://publications.sto.nato.int/publications/STO%20Meeting%20Proceedings/STO-MP-HFM-361/MP-HFM-361-P13.pdf>
- Nestor, M.W. and Wilson, R.L. (2022). *Anticipatory ethics and the use of CRISPR in humans*. Springer.
- Nissenbaum, H. (2009). *Privacy in Context: Technology, Policy, and the Integrity of Social Life*. Stanford University Press.
- O'Neil, C. (2016). *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. Crown.
- Platt, M., Platt, D., & McBurney, P. (2024). Sybil attack vulnerability trilemma. *International Journal of Parallel, Emergent and Distributed Systems*, 39(3), 446–460. <https://doi.org/10.1080/17445760.2024.2352740>
- Pariser, E. (2011). *The Filter Bubble: What the Internet Is Hiding from You*. Penguin Press.
- Singer, P. W., & Brooking, E. T. (2018). *LikeWar: The Weaponization of Social Media*. Eamon Dolan/Houghton Mifflin Harcourt.
- Sociotechnical Imperative. [Sociotechnical Imperative → Area → Resource 5](#). Accessed 2/5/2026.
- Sunstein, C. R. (2017). *#Republic: Divided Democracy in the Age of Social Media*. Princeton University Press.
- Tavani, H.T. (2009). *Ethics and Technology: Ethical Issues in an Age of Information and Communication Technology*. Wiley.
- Tomassi, A., Falegnami, A., & Romano, E. (2024). Mapping automatic social media information disorder: The role of bots and AI in spreading misleading information in society. *PLOS ONE*, 19(5), e0303183. <https://doi.org/10.1371/journal.pone.0303183>

- US Cyber Command. (2024, September 3). *Russian Disinformation Campaign "DoppelGänger" Unmasked*.
<https://www.cybercom.mil/Media/News/Article/3895345/russian-disinformation-campaign-doppelgnger-unmasked-a-web-of-deception/>
- Vosoughi, S., Roy, D., & Aral, S. (2018). The spread of true and false news online. *Science*, 359(6380), 1146-1151.
- Wilson, Richard. Anticipatory Ethics as a Method for Teaching Engineering Ethics. 2021 ASEE St. Lawrence Section Conference.
- Wilson, Richard L. Cambridge Analytica, Facebook, and Influence Operations: A Case Study and Anticipatory Ethical Analysis, Proceedings of the 18th European Conference on Cyber Warfare and Security, University of Coimbra, Portugal, Academic Conference and Publishing International, Ltd. 2019.
- Wilson, Richard L., Cyber Warfare, Terrorist Narratives and Counter Terrorist Narratives: An Anticipatory Ethical Analysis, Proceedings of the 18th European Conference on Cyber Warfare and Security, University of Coimbra, Portugal, Academic Conference and Publishing International, Ltd. 2019.
- Wilson, Richard L. Information Warfare: Fabrication, Distortion and Disinformation: A Case Study and Anticipatory Ethical Analysis, Proceedings of the 18th European Conference on Cyber Warfare and Security, University of Coimbra, Portugal, Academic Conference and Publishing International, Ltd. 2019.
- Zuboff, S. (2019). *The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power*. PublicAffairs.